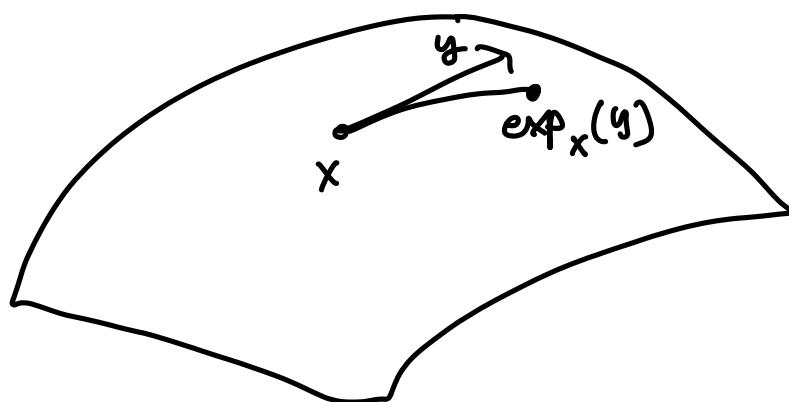


Query complexity of Riemannian optimization

- Outline :
- 1) Quick Riemann opt. refresher
 - 2) Examples; hyperbolic plane
 - 3) Upper bounds
 - 4) Hamilton-Moitra lower bound
 - 5) Open directions

Riemannian Opt : M Riemannian manifold

- Metric $\|\cdot\|_x$
- $d(x, x') := \text{length of shortest } x \rightarrow x' \text{ curve}$
= length of geodesic between.



- assume we have access to the exponential map $\exp_x : T_x M \rightarrow M$; goes distance $\|y\|_x$ along geodesic in direction y .
- gradient $\nabla f(x) := \nabla_{y=0} f \circ \exp_x(y)$
"direction of steepest geodesic".
- First order optimization:
query at x , get $f(x), \nabla f(x)$
- want to find point x s.t.
 $\|\nabla_x f\| < \varepsilon$, $|f(x) - \text{OPT}| < \varepsilon$, or $|\text{d}(x, x^*)| < \varepsilon$.

- convex: $\partial_t^2 f \circ \exp_x(yt) \geq 0 \quad \forall y, \forall x \in M$
- α -strongly convex: $\partial_t^2 f \circ \exp_x(yt) \geq \alpha \|y\|_x^2$.
- β -smooth: $\partial_t^2 f \circ \exp_x(yt) \leq \beta \|y\|_x^2$

Manifold examples:

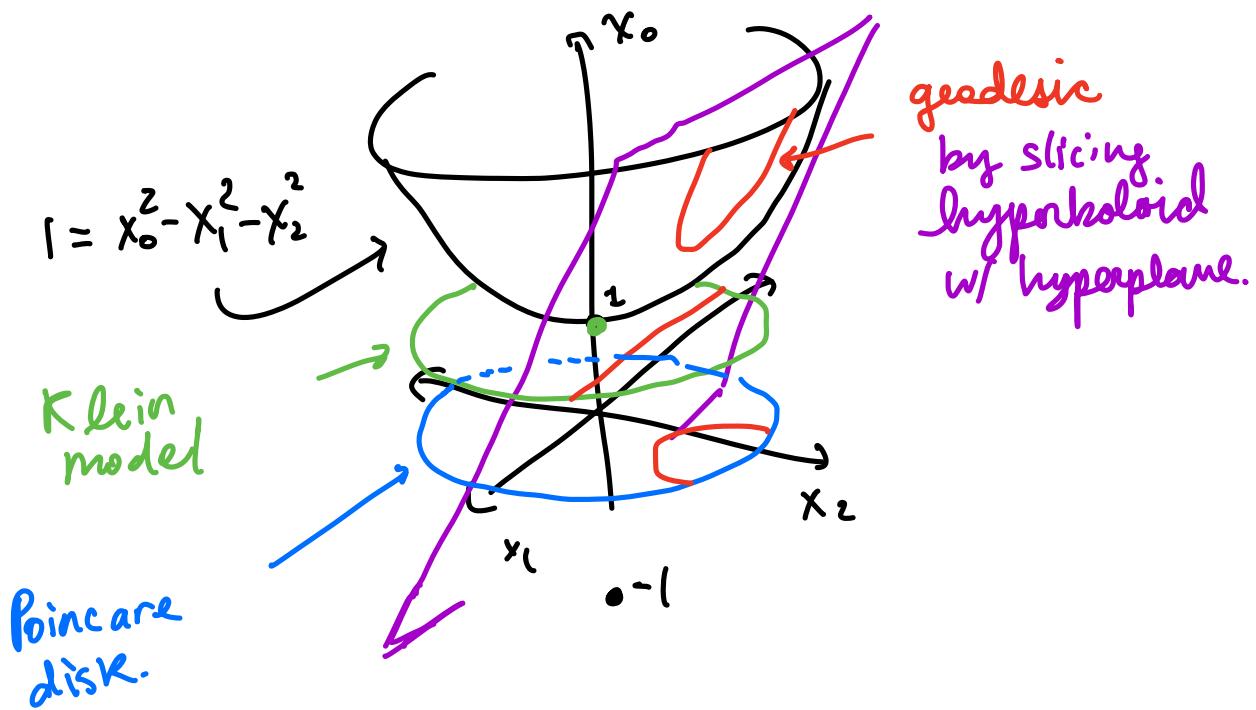
- $\text{PD}(n)$, $\exp_x(Y) = \sqrt{x} e^Y \sqrt{x}$. \mathbb{R} or \mathbb{C}

- $\text{PD}_1(n) = \text{PD}(n) \cap \text{SL}(n)$
- Hyperbolic plane: $\text{PD}_1(n, \mathbb{R})$.

$$\det \begin{bmatrix} x_0 + x_1 & x_2 \\ x_2 & x_0 - x_1 \end{bmatrix} = 1, \Leftrightarrow x_0^2 - x_1^2 - x_2^2 = 1$$

$x_0 > 0$

"hyperboloid
model"



Function examples

- Kronecker covariance estimation:

Given tensor $v \in \mathbb{R}^{d^3}$, minimize

$$\|(x_1 \otimes x_2 \otimes x_3) v\| \text{ over } \text{PD}(d)^3.$$

AKA "tensor scaling".

- Matrix means: given P_1, \dots, P_n in $\text{PD}(d)$,
minimize $\sum_i d(x, P_i)^2$ over $\text{PD}(d)$.

Examples so far are

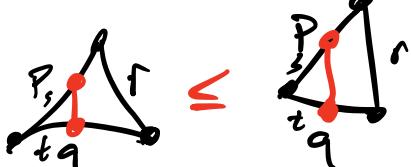
Hadamard manifolds, i.e. Riem. manifolds that:

- Have sectional curvature ≤ 0 &
- Are simply connected.

→ \exp_x of 2-D subspaces look like pringles

Properties: if M Hadamard, then:

- Triangle comparison:



- Hadamard-Cartan: $\exp_x: \mathbb{R}^n \rightarrow M$ diffeomorphism.

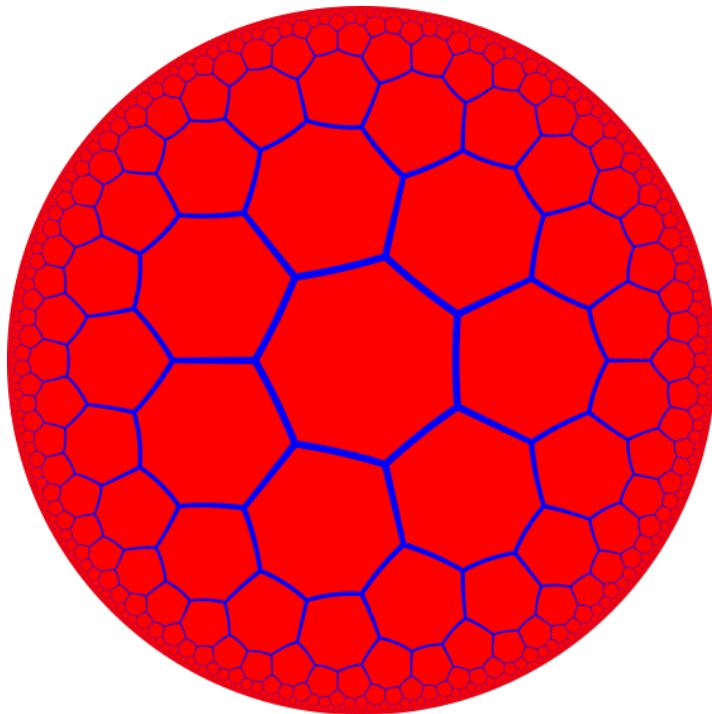
- Bishop-Gromov:

$$\text{vol}(B_\varepsilon \subseteq M) \geq \text{vol}(B_\varepsilon \subseteq \mathbb{R}^n)$$

- Distance-convexity:

$d(x, \cdot)^2$ strictly convex.

Key fact about H_2 : $\text{vol } B_x(r) = e^{\Theta(r)}$.



Poincaré disk model.

Positive algorithmic results:

Parameters: β smoothness, K condition # $\frac{q}{\beta}$
 R , radius of ball \geq domain.
 $T = \# \text{ steps}$, n dimension.

Class	Algorithm	Guarantee
β -smooth	GD	$\ \nabla f\ \leq \sqrt{\frac{\beta(f(x_0) - f(x^*))}{T}}$
	Folklore	$f(x) - f(x^*) \leq \frac{\beta R^2}{T}$
K -conditioned	AGD Martinez-Rubio 2020	$f(x) - f(x^*) \leq \frac{\beta}{T^2} e^{O(R)} *$ vs $\frac{LR^2}{T^2}$ Euc.
	GD Zhang, Sra → 2016	$\ \nabla f(x)\ \leq \ \nabla f(x_0)\ e^{-T/K}$ $f(x) - f(x^*) \leq \beta R^2 e^{-T/K}$
	AGD Martinez-Rubio 2020	$f(x) - f(x^*) \leq \beta e^R e^{-T/\sqrt{R}} *$ vs $LR^2 e^{-T/\sqrt{R}}$ Euc.

* Only on H_n .

L lipschitz,
differentiable

Cutting Planes	$f(x) - f(x^*) \leq n L e^{-T/n^2}$
Rusciand '18	$\underbrace{n R}_{\text{volume of ball}} e^{-T/n^2}$
Non-algorithmic!	Euc. $L R e^{-T}$.

Issues: Even in H_2 , to get

- AGD needs $R \sqrt{K}$ queries
vs $\sqrt{K} \log R$ Euclidean
- Cutting planes needs R
vs $\log R$ Euclidean.

Question: Do Euclidean bounds
carry over??

Answer: Probably not.

Work of Hamilton - Moitra 2020 suggests:

Conjecture:

- For 1 lipschitz, need $\Omega(R)$.
- For $O(R)$ conditioned, need $\Omega(R)$.
(not $\sqrt{R} \log R$)

Comment: $O(R)$ is best possible condition # on $B(R) \subseteq H_2$. (HM13).

$\sqrt{K} \sqrt{R}$?

- They prove this in a noisy model.
- Noise can be R^{-100} ; wouldn't cause problems in \mathbb{R}^n .
- Why do I care? In tensor scaling problem, need $R = \exp(n)$!

Noisy Query Model :

- Unknown differentiable function f w/ minimum x^* in $B(R)$.
- Can query any point $x \in B_{100R}(0)$, get $f(x) + z_0$, $\nabla f(x) + z_1$, where z_0, z_1 are independent (& do not depend on past queries).
- Say z (noise) is c -nonconcentrated if the density function of z is everywhere $\leq c$, δ -precise if $\|z\| \leq \delta$ always.

- Concretely: let $z_2 \sim \text{unif } B_{\mathbb{R}^2}(\delta)$; $c = \frac{1}{\pi \delta^2}$
 $z_1 = \text{unif } B_{\mathbb{R}}(\delta)$; $c = \frac{1}{2\delta}$.

KEY FACT: Small (precise) noise, doesn't interfere w/ Euclidean optimization.

Main Theorem There are function classes

- F_c $O(R)$ -conditioned
- F_L $O(1)$ -Lipschitz

for which finding a point $\leq \frac{R}{5}$ from x^*

requires

$$\Omega\left(\frac{R}{\log R + \log \delta + \log c}\right)$$

δ -precise, c -nonconvex noisy queries.

e.g. $\delta = R^{-10}$, $c = R^{20}$. $\rightsquigarrow \Omega\left(\frac{R}{\log R}\right)$.

Remarks:

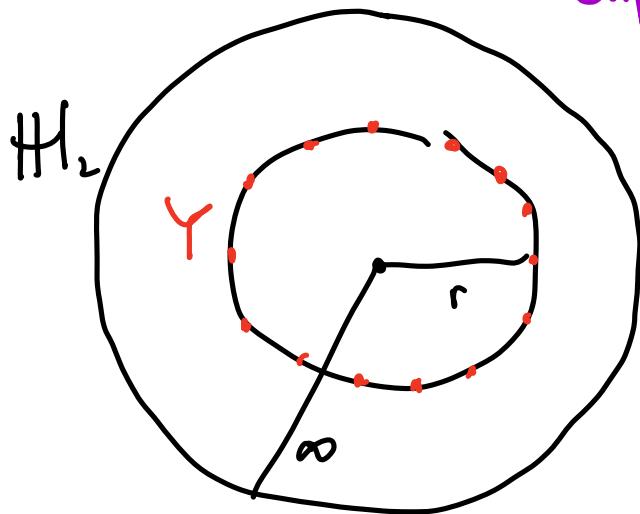
- Noise doesn't preclude $\log(\frac{1}{\alpha})$,
just $\log(R)$, \sqrt{R} dependence.
- can replace $d(x, x^*) \leq \frac{R}{S}$ by
 $f(x) - f(x^*) \leq \frac{R^2}{S}$ or $\|\nabla f(x)\| \leq \frac{R}{S}$.
- many classes would work.
- also works for $\nabla^2 f$ queries.

Proof of HM'21

- function class: class consists of
functions f_y , $f_y(x) = \hat{d}(x, y)$. (for $k \leq r$),

or $f_y(x) = d(x, y)^*$ (for lipschitz).

- Noise necessary: else gradient points exactly towards $y = \operatorname{argmin} f_y(x)$!
- Which y ? y is chosen from a Known set \mathcal{Y} of $n = e^{\Theta(R)}$ equally spaced points at distance $\frac{R}{2}$ around unit circle in H_2 . Can fit b/c vol exponential in R !



Figuring out which element $y \in \mathcal{Y} \Leftrightarrow$
Finding point within $\frac{R}{4}$ of $\operatorname{argmin} f$.

Game formulation:

- Alice knows an element $y \in \mathcal{Y}, |\mathcal{Y}|=n$.
- Bob makes a query $q \in Q := \mathcal{B}_{1000R}(\delta)$

receives observation $x \in \mathcal{X}$,

Sample from distribution

$p_{q,y}$:= density of pair $(f(q) + z_1, \nabla f_y(q) + \vec{z}_2)$
↑
noise.

here $x \in \mathcal{X} := [0, (1000R + \delta)^2] \times \mathcal{B}_{R^2}(1000R + \delta)$

- Bob wins if he guesses y w/
probability $\geq 2/3$.

[HM13] proves a general lower
bound on noisy query games:

Theorem: For any noisy query game

$$(\mathcal{Y}, \mathcal{Q}, \mathcal{X}, P)$$

Alice's secrets \uparrow queries \uparrow observations \uparrow map $p: \mathcal{Q} \times \mathcal{Y} \rightarrow P(\mathcal{X})$

such that $P_{x,y}$ is c -nonconcentrated,
 then Bob needs $\Omega\left(\frac{\log |\mathcal{Y}|}{\log(c \text{ vol}(X))}\right)$
 queries to win.

Corollary: For noisy gradient task,

$$|\mathcal{Y}| = e^{\Theta(R)}, \quad \text{vol}(X) = O((r^2 + \delta)(r + \delta)^2).$$

$$\Rightarrow \text{vol}(X) = O(r^3 \delta^3),$$

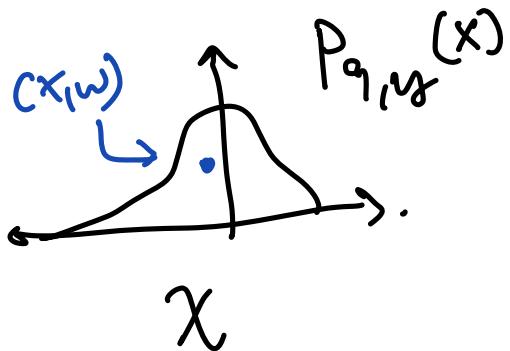
$\Rightarrow \Omega\left(\frac{R}{\log R + \log \delta + \log c}\right)$ queries
 necessary.

Proof of NQG theorem:

- prove lower bound for game where Bob gets even more info.

Def: Transparent NQG: Same as

NQG, except Bob gets sample from uniformly random point (x, w) under graph of $P_{q,y}: X \rightarrow \mathbb{R}$



marginal distrib of x is $P_{q,y}$;
only gives Bob more info.

If Alice samples $y \sim \text{Unif}(Y)$, Bob's optimal strategy after t steps is to output y^* most likely under posterior

$$y | q_1, (x_1, w_1), q_2, \dots, q_{t-1}, (x_t, w_t)$$

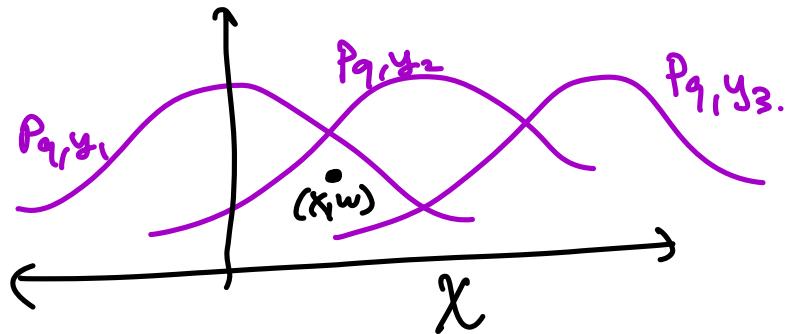
all outcomes so far

Lemma: Let $y \sim \text{Unif}(Y)$, consider TNQG.

- 1) posterior for y after 1 query
is uniform on some $\mathcal{Y}_1 \subseteq Y$
- 2) in expectation over observation,
 $\log |\mathcal{Y}_1| \geq \log |\mathcal{Y}| - \log(\text{Vol}(X))$.

Proof: 1) After query $q_1, (x, w)$
sampled at uniform from

graph of $P_{q,y}$ for $y \sim \text{Unif}(\mathcal{Y})$.



Claim: $\mathcal{Y}_1 = \{y' : (x, \omega) \in \text{Graph } P_{x, y'}\}$.

- After query, Bob knows $y \neq y'$ for any $y' \notin \mathcal{Y}_1$.
- Moreover, because every graph has area 1, posterior is uniform on all $y' \in \mathcal{Y}_1$.

2) Let $N(x, \omega) := |\mathcal{Y}_1|$. Then

$$\mathbb{E}_{y} \left[\int_{\text{Graph } P_{q,y}} \log N(x, \omega) dx dw \right]$$

$$= \mathbb{E}_{\text{noncon.}} \left[\int_{X \times [0, C]} \underset{\text{Graph } P_{q,w}}{1} \log N(x, w) dx dw \right]$$

$$= \int_{X \times [0, C]} \frac{N(x, w)}{|Y_1|} \log N(x, w) dx dw *$$

► Each graph has area 1 $\Rightarrow \int_{X \times [0, C]} \frac{N(x, w)}{|Y_1|} dx dw = 1.$

► Thus * is - entropy of $N(x, w)/|Y_1|$ up to additive const;

$\Rightarrow *$ minimized

when $N(x, w) \equiv \text{const.}$ over $X \times [0, C] \Leftrightarrow \frac{|Y_1|}{C \text{vol}(X)}$.

► Min value for *:

$$\int \frac{1}{C \text{vol}(X)} \log \left(\frac{|Y_1|}{C \text{vol}(X)} \right) dx dw$$

$$= \log \left(\frac{|Y_1|}{C \text{vol}(X)} \right) = \log |Y_1| - \log(C \text{vol}(X)).$$



- Now using Lemma, prove NQG theorem.

- By induction, after t queries, posterior uniform on Υ_t with

$$\mathbb{E} \log |\gamma_t| \geq \log |\gamma| - T \log [c \text{vol}(X)]$$

or $Z := \log |\gamma| - \log |\gamma_t|$

has $\mathbb{E} Z \leq T \log [c \text{vol}(X)]$.

Algorithm only succeeds if $|\gamma_t| = 1$,

i.e. $Z \geq \log |\gamma|$.

By Markov, $\Pr[Z] \geq T \log [c \text{vol}(X)]$

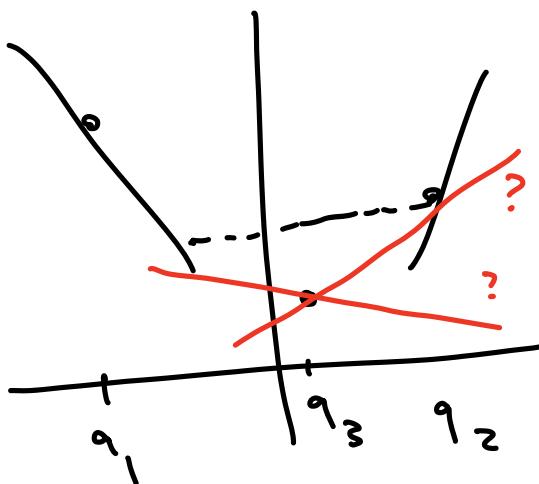
$\leq \frac{1}{3}$; hence win with

$$\Pr \leq \frac{1}{3} \text{ if } T \leq \frac{\log |\gamma|}{3 \log [c \text{vol}(X)]}. \quad \square$$

Open directions:

- Removing noise?

Play "20 questions": $\frac{1}{2}$ search space



Issue: which gradients consistent
w/ previous observations?

Nontrivial in H_2 ; no "linear functions"
 $\nabla \cdot \log_x(y)$ not convex in y in

- Algorithmic cutting planes/
ellipsoid in $H_n \dots H_2$ maybe O.K.

- Function classes that can be optimized in $O(\log R)$

queries? & include interesting problems like tensor scaling 😊